

Original Article

Optimizing Cloud-Native Development Costs: Smart Spending in the Cloud

Rakesh Kumar Mali

Delivery Module Lead, Atlanta, Georgia, USA.

Corresponding Author : rakesh.mali.jmd@gmail.com

Received: 06 June 2025

Revised: 28 June 2025

Accepted: 19 July 2025

Published: 30 July 2025

Abstract - Effective cost management has become more crucial as cloud adoption increases across all sectors. For Technical Program Managers (TPMs) in charge of software development projects, management has become a crucial skill. This work offers a thorough framework for cloud cost optimization, combining financial management with technical best practices. And examine essential tactics like granular cost transparency, resource rightsizing, policy implementation, and performance-cost balancing. Real-world case studies and empirical data support the conclusion that TPMs can deliver substantial cost savings without compromising innovation speed or service quality. This work presents a maturity model for cloud financial operations (FinOps) and discusses areas for future investigation and guidelines for this ever-changing sector.

Keywords - Cloud Computing, Software Engineering Economics, Technical Program Management, Cost Optimization, FinOps.

1. Introduction

Cloud-native architectures have completely changed the economics of software development and deployment. Once restricted by the capital costs of on-premises infrastructures, organizations can now opt for elastic consumption models that charge only for actual usage. Worldwide end-user spending on public cloud services is forecast to grow 23.1% in 2021 to total \$332.3 billion, up from \$270 billion in 2020. This meteoric expansion highlights the rapid uptake of cloud technology throughout a range of industries, prompted largely by requirements for scale and availability combined with speed to market.

But the speed at which organizations can create resources can outpace their ability to manage and optimize costs. When left unchecked, cloud costs can spiral out of control, driven by a combination of over-provisioning, idle workloads, and poor architectural decisions that lead to inefficiency. As reported by Flexera's State of the Cloud 2022, enterprises waste 32% of their cloud spend — a shocking number that illustrates the need for disciplined cost management in the rapidly growing world of cloud computing. To do this, Technical Program Managers play a critical role in connecting engineering teams focused on feature delivery with finance teams focused on cost. However, their job is not easy, as they need to cooperate with different priorities; on the one hand, development should not be held back by high cost constraints, and on the other, the cloud investment should not run out of control. Thus, effective TPMs apply the principles of FinOps, bridging the gap between DevOps, finance, business leadership, and cloud investment objectives.

In this article, we compile best practices and pioneering techniques so TPMs can take actions to ensure efficient cloud cost while not sacrificing innovation or quality. From understanding how to 'rightsize' large workloads and set up automated scaling policies, right through to negotiating committed-use discounts and building a cost-aware engineering culture: We will cover practical steps you can take towards reducing cloud spend without sacrificing performance or scalability.

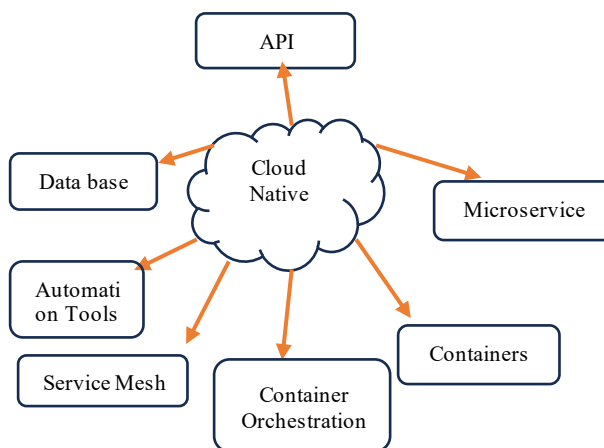


Fig. 1 Cloud native architecture

By embracing the power of the cloud to its fullest potential, i.e., scalability, automation, resilience, and cloud-native architecture, the building, deployment, and management of software can be taken to a new level. Unlike traditional monolithic systems that exist on physical servers, cloud-native applications are created as loosely coupled microservices; a collection of specialized



units that run in lightweight containers (e.g. Docker) that are managed using a platform, such as Kubernetes, which is used to automate deployment, scaling, and recovery. This modular design feature has the advantage of letting teams make changes to individual components at an increased pace of agility, but it is also accompanied by complexity in cost tracking. The added efficiencies include the use of serverless codes (e.g., AWS Lambda), where codes only run when triggered, cutting idle resource costs, and event-driven architectures (e.g., Kafka), decoupled services to eliminate faults. Yet, the same resource elasticity provides the foundation for innovation, the auto-scalable nature of resources, and a highly optimized CI/CD pipeline can cause cost blow-outs unless properly managed. Such things can be a very overprovisioned Kubernetes cluster, orphaned cloud storage, or resources with no tags, and enterprises end up wasting an average of 32% of their cloud expenditure (Flexera 2022). TPMs should fill the gap by incorporating FinOps methods into the dev processes: real-time cost dashboards (e.g., AWS Cost Explorer), tagging policies to assign costs by team or project, and rightsizing resources based on actual metrics of usage. They also make strategic tradeoffs, i.e., putting workloads on spot instances to make significant cost savings (up to 90 percent) or reserved instances to have more predictable needs in a stable production system. That said, cloud-native is not just a technical concept since it is also a cultural paradigm wherein engineering, financial, and operations teams are able to work together to balance cloud spending and business results so scalability does not have to come at the expense of financial discipline.

Cloud-native development has transformed the bespoke construction of scalable applications through containerization and microservices, DevOps, and continuous integration/delivery (CI/CD) pipelines. With more businesses moving to the cloud, the sleek factor of enhanced agility, scale, and quick deployment is sometimes offset by the looming shadow of ever-rising operational costs. In most cases, over-provisioning of resources, poor usage of services, over-provisioning of compute-storage, and non-optimization of resources occur when cloud is adopted without planning costs.

1.1. Problem Statement

Although most companies accept the use of cloud-native architectures, they still experience severe challenges regarding cost reimbursement for cloud equipment and optimization. Recent cloud cost optimization tools are inclined to an approach that would position the cost intelligence as a post-deployment analytics process instead of integrating it into the application design and at runtime.

This reactive model does not present real-time suggestions or dynamically changing strategies, keeping pace with the organization's working patterns and cost limitations. Hence, an urgent need emerges on the one hand to promulgate an active, smart, and systematic mechanism that incorporates the idea of cost-sensitivity into the cloud-native application development process.

1.2. Research Gap

Several works have been done on cost optimization of clouds at the reserved instance level, autoscaling, and price model level. Nevertheless, there is the dominance of vendor-specific solutions, solutions that are not generalizable across hybrid and multi-cloud environments and solutions that are purely static in nature. In addition, although AI and ML models have been adopted to predict cloud workloads and allocate resources, only a few frameworks have proposed cost-efficiency as a dynamic objective alongside performance indicators. The literature on cost optimization and cloud-native architectural patterns leads to little addressing of a model that shows the internal existence of both at once, i.e., the unified, adaptive, and generalizable model.

This article fills such a research gap by suggesting the SISSGECO (Smart Integrated Strategy of Scalable Generalized Cloud Optimization) model. The model dwells upon real-time cost-consciousness, workload flexibility, and scalable deployment plans. The framework presents a unique performance-cost optimization mechanism, which is accelerated through smart feature choice and resource allocation prioritization, focused on addressing inefficiencies witnessed in the current cloud-native method of development.

1.3. Research Gap with Existing Work

Current approaches to optimization of cloud costs are all either static or platform-specific, or only at the scale of infrastructure. They do not tend to be flexible, work in real time, and be agile with cloud-native developmental processes. The existing construction methods seldom integrate high-end feature choice with hybrid machine learning to create dynamic choices of cost-performance.

As seen in the literature, a distinct lack of a model is actively used to target the current problem on a unified, ML-driven approach of general cost optimization across the CI/CD pipeline in the multi-cloud setting. Intelligent feature selection and hybrid classification work together to help solve this issue, as proposed in the SISSGECO model to align cost decisions with real-time dependencies, workloads and scale.

2. Literature Review

The problem of cloud cost optimization has developed into a serious issue in contemporary enterprise architectures, with the rapid adoption of cloud-native designs becoming popular all over the business environment [1,2]. The initial studies mostly rested on reactive cost-saving methods, whereby resource provisioning was altered through rules based on thresholds and scheduled use patterns [3]. Although successful when applied to constant workload, they failed in variable-demand and mixed workloads with increasing and decreasing demand [4,5].

Later experiments presented the idea of autoscaling, where resource allocations change in accordance with the

current parameters, like CPU idle time or memory usage [6,7]. These methods enhanced responsiveness but are often clueless, leading to over-provisioning or performance impediments [8]. Also, most of the solutions were designed to be used in one platform or in one provider and did not fit well in multi-cloud or hybrid cloud environments [9].

This and the predictive approaches to modelling demand and distribution of resources initially began to attract a hearty response when machine learning methods caught on [10,11]. Some of the initial techniques involved in estimating resource requirements and optimizing cost on the infrastructure included linear regression, support vector machines and decision trees [12]. Nonetheless, these models proved incapable of generalizing to different types of applications and are intolerant of input variability.

Most recently, deep learning and reinforcement learning have been used to train smarter autoscaling and scheduling, to realize better cost-performance tradeoffs [14,15]. Others have added an uncertainty-aware prediction and feedback optimization loop [16]. These achievements notwithstanding, there is still a gap in implementing an efficient hybridization of intelligent feature selection and flexible and cost-effective classification models [17,18].

The existing approaches mostly have rather isolated layers, whether infrastructure, orchestration or application, without the integration of cost intelligence into the development lifecycle [19]. The entrenchment of optimization activities in CI/CD pipelines and releases through microservices is also not given as much attention. This perceived fragmentation has brought about an urgent demand to develop a single, scalable, flexible framework that strikes a balance between cost and performance in real time and on a request basis, at least in a world in which organizations run in fast-changing cloud-native landscapes.

3. Cloud Cost Visibility and Attribution

Cloud-native cost visibility and attribution provides financial governance with a foundation in cloud cost and spend management, as cloud spend and cloud systems and organizations track, understand and optimize cloud spend. Without granular visibility, enterprises can end up with uncontrolled spending related to shadow IT, overprovisioned resources, and so on, which is even more problematic to resolve when dealing with auto-scaling microservices and ephemeral server-based workloads. The tagging approach (e.g., labeling the resources by the project/reports, department, or the environment) and the use of the hierarchical account structures (e.g., AWS Organizations, Azure Management Groups) are necessary to attribute the costs to the business units. Tooling such as AWS Cost Explorer, Azure Cost Management, or 3rd party solutions such as CloudHealth will include real-time dashboards, anomaly detection, and forecasting, all of which are heavily dependent on consistent metadata (e.g.

tagging compliance). However, there are multi-cloud/hybrid cases where there are limitations, in that the different billing systems complicate aggregation of reporting and how it should charge on shared items (e.g., Kubernetes cluster, databases) either on a per-namespace or per-usage basis. According to Gartner (2023), firms with well-developed cost-visibility processes save money on the cloud by 25 to 40 percent, whereas those that skip an attribution process experience the phenomenon of bill shock, when little-tracked experimental workloads driven by experimentation cause shocks when they appear in bills. An effective implementation of accountability implies cross-functional cooperation: engineering teams should focus on cost-efficient production (e.g., use of cost-efficient types of instances), the finance team should establish its budgetary limits with warning signs, and the leadership must synchronize cloud investments with ROI achievement. Innovation opportunities in areas of AI-enhanced cost anomaly detection and auto-enforcement of such policies (i.e.: shutting down non-prod resources outside regular business hours) are also emergent; however, the cultural acceptance of increased transparency of the costs (and thus use) of various resources remains one of the biggest impediments in arriving at a fully accurate attribution of costs back to their owners and consumers.

3.1. The Challenge of Distributed Architectures

Among the most basic complexities of controlling cloud spend is the ability to clearly understand expenditure across distributed systems. The cloud native environments, especially on microservice-based architecture, distribute the spending into hundreds or even thousands of dynamically sized services, containers and serverless functions, unlike monolithic architecture, where the spending is centralized and therefore easier to track. Although the benefits microservices bring in terms of scalability, fault isolation and maintainability are indisputable, it is inherent in their structure that they make the exact cost of the particular feature or a part of the product practically unknowable. For example, there can be only one user-facing feature comprising multiple APIs, databases, and event-driven workflows that exist in various containers and in different availability zones with costs buried under layers of infrastructure abstraction. This cost dispersion is compounded by auto-scaling policies, spot instances and cross-service relations, and it is not always easy to even answer basic questions such as, “How much does the checkout service cost per transaction?” or “Whose team gets this increase in S3 storage costs?”

The seriousness of this difficulty is confirmed by industry research: 37% of respondents in Flexera 2023 Cloud Report have listed the lack of visibility as the most significant cost management barrier in the cloud, and Gartner has also reported that an average of 30% of total budgets is cloud spending that is not attributed. This is escalated by the fact that the problem is compounded by the fact that:

3.1.1. Shared Resource Contention

Kubernetes clusters, message queues (e.g., Kafka), and data lakes (e.g., Snowflake) will usually be shared, so a cost allocation model that may be proportional (e.g., by CPU-hours or memory usage per namespace) will be required. Without uniform changes in tagging or such tools as OpenCost, costs become opaque, and accountability is lost.

3.1.2. Ephemeral and Serverless Workloads

Small-life containers (e.g., CI/CD pipelines) or serverless functions (e.g., AWS Lambda) introduce the concept of short-lived costs as one-off (transient) and not easily attributable, even using traditional monitoring tools. Sub-second billing of AWS Lambda, for example, can support millions of micro-transactions a month that do not help identify cost drivers.

3.1.3. Multi-Tenant and Multi-Cloud Complexity

Inconsistent labeling, currency being out of sync, and different month-to-month billing strategies (AWS CUR and Azure Cost management) will wind up breaking a cross-board sight. As reported by IDC in 2023, just 17% of enterprises identified that they possess unified multi-cloud cost reporting.

3.2. Implementing Robust Tagging Strategies

Technical Program Managers become instrumental to the pursuit of Cloud Cost Accountability because they help coordinate activities across departments, such as finance and technical teams. The first and most important step is to introduce sound cost allocation tagging systems, which will allow precise mapping of expenditures to business units, projects, deployment environments and even feature-level usage. The following three aspects are the key factors to be addressed in its effective implementation:

- The common taxonomy of tagging and naming is applied on an organizational-wide scale.
- The uniform type of infrastructures-as-a-code templates and pipeline injection
- Regular checkups of governance to keep track of surfaces close to business arrangements

3.3. Leveraging Cost Management Platforms

Newer cloud cost intelligence tools, such as CloudZero and Vantage, enable raw cloud billing data to be converted into useful actions through the power of cloud billing analytics. In these solutions, data about the spent money summarizes across different cloud providers and provides configurable visualization dashboards to visualize the spending patterns. In addition to the basic reporting abilities, they use machine learning algorithms to detect anomalous real-time cost fluctuations. The most important financial and operational indicators that organizations need to monitor are as follows:

- YoY and MoM of expenditure analysis
- Unit economic measures (cost/customer/transaction)

- Compute/storage consumption benchmarks of resources
- Solution of cost in terms of service level (AWS/Azure services categories)

3.4. Case Study: A Fintech's Cost Attribution Journey

A financial technology firm set up a complete cost formation system using CloudZero. By cost allocating their billing to internal cost centres and teams, they gained a fine delved into the spending activities. Nonetheless, difficulties were presented:

- Manual mistake of tagging and omission
- Problems in assigning joint resources, such as Reserved Instances
- Marketplace purchase attributions can be delayed

The Fintech to meet these problems has adopted:

- Tag enforcement at the level of the cloud platform based on policy
- They used it as part of their infrastructure-as-code solution to automatically tag
- A special script will assign the cost of the Reservation Instance based on actual usage
- A new process for real-time attribution of purchases in marketplaces

These profits meant a 15 percent reduction in uncredited spending and enabled us to estimate the team members better and hold them accountable.

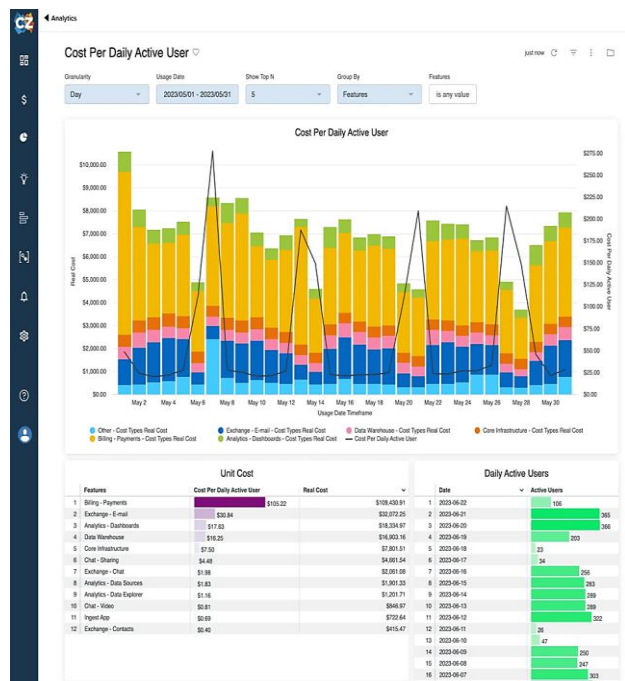


Fig 2 CloudZero Cost Dashboard

4. Rightsizing and Optimizing Cloud Resources

Rightsizing is also one of the most valuable cloud cost control techniques; it means aligning the amount of cloud resources with the real workloads. A lot of entities end up overspending when they over-provision, just to be safe; they end up with underutilized virtual machines, bloated

storage, and idle databases. The waste that rightsizing eliminates can be characterized by the analysis of performance indicators such as CPU utilization, memory utilization, consumption, and network throughput to address instance type and storage tier to the actual demand related thereto. As an illustration, an AWS EC2 instance that runs at 15 percent CPU utilization can perhaps be scaled down to a smaller instance or even switch to serverless implementations such as AWS Lambda that may handle the occasional jobs without affecting performance.

When it comes to comprehensive cloud optimization, though, it needs a multi-dimensional approach beyond mere downsizing. Spot instances offer the opportunity to achieve a compute cost reduction up to 90 percent compared to sustained workloads, and auto-scaling policies keep the resources scaled up and down according to the traffic patterns. Further waste optimization can be fulfilled by optimizing storage, e.g., storing more used data in cheaper storage tiers (e.g. AWS S3 Glacier) or by using life cycle data policies. Newer technologies under FinOps, like making Rightsizing recommendations on the AWS Cost Explorer or using third-party tools like Densify, train in-house machine learning models, which make determinations about allocation based on historical consumption and the unique actions that are then recommended to an organization, eliminating guesswork.

Even so, rightsizing is not an occasional action, and it requires being monitored and altered. Cloud environments change at high rates, application usage rising or dropping on seasonal grounds, on feature release or on behavioural changes of the user. A systematic review process must be instituted where teams can review the utilization metrics on a regular basis and make changes in resource allocation to save on costs in the long run. By rightsizing the right way, you not only save money but also increase performance levels through the removal of resource contention and lessen environmental impact by using energy more efficiently. Automated tooling and cross-team collaboration allow organizations to find the balance between operational readiness and cost efficiency in their cloud environment, which can be subject to very fragile equations.

4.1. Embracing Cloud-Native Architectures

The Technical Program Managers at the organization must lead the effort to adopt cloud-native solutions, especially containerization (e.g., Kubernetes) and serverless compute (e.g., AWS Lambda), to get dynamic demand allocation of resources. Through these newer paradigm architectures, infrastructure can scale to align perfectly with the workload requirements, and wastage associated with over-provisioning can be avoided in the older systems. Real-world experience proves the practical advantages: The companies that move monolithic apps to microservice containers with containerization regularly save 20-50 percent in terms of efficiency cost optimizations. In addition to the actual savings, these

architectures save deployment and scaling overheads through automation of such tasks.

4.2. Identifying and Eliminating Waste

Some of the most common causes of inefficiency in clouds include over-provisioning of infrastructure and wasted computing capacity. Technical program managers are doing exceptionally important work because they are joining forces with infrastructure teams to bring intelligent rightsizing solutions. Such automated systems constantly monitor resource utilization statistics and provide data-assisted suggestions of an optimum instance configuration such that some desired workload can be run on the infrastructure of the right size.

4.3. Leveraging AI-Driven Optimization

Innovative services such as ProsperOps use machine learning models to constantly examine the usage behavior and automatically fine-tune the Reserved Instance/Savings Plans portfolio. This smart automation consistently performs 15-40 percent better cost savings than a non-dynamic manually handled discount unit.

4.4. Automated Instance Purchasing

The intelligent reservation management platforms that can be used, including ProsperOps, can provide significant cloud cost optimization. These AI-driven technologies continuously analyze workload tendencies and automatically purchase the most beneficial mix of Reserved Instances and Savings Plans with maximum discounts and the least risks of commitment obligations.

4.5. Case Study: Ephemeral Environment Optimization

A recent ephemeral environment optimization project provided a revolutionary outcome: By using intelligent cluster rightsizing of Kubernetes clusters, it also saved 48.1 percentage points off an otherwise high idle resource cost of 71.3%, bringing the figure to 37, a relative improvement of 48%. The fiscal savings were further augmented by surprise performance increases, with the most influential workload output times speeding up 3-4 hours per job run, resulting in a scripted time of acceleration of workload processing time accordingly.

5. Enforcing Cost Governance

The solution to effective cloud cost governance is to have clear policies along with automated guardrails and create accountability to help avoid overrun on the budgets, but not degrade the agility of operations. This comes in the form of setting limits and warnings on budgets using tools such as AWS Budgets or Azure Cost Management, enforcing consistency with tagging to allow proper cost attribution, and creating approvals on high-cost services.

Creating a pipeline in which governance checks are built into CI/CD (e.g., policy-as-code with Open Policy Agent of Terraform) and regular finance-oriented reviews will allow organizations to balance financial control and organizational innovation needs.

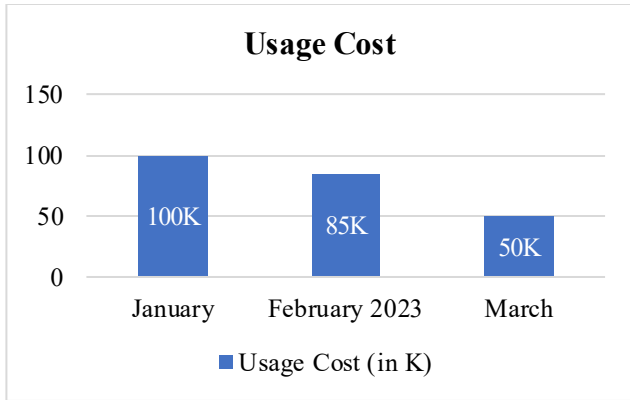


Fig. 3 Azure Monthly Usage Cost Dashboard

Monthly Usage and Cost by Service						
Month	February		January		March	
MeterCategory	No. of Units	Cost	No. of Units	Cost	No. of Units	Cost
Azure Cosmos DB	20000	\$20,000	30000	\$40,000	10000	\$10,000
SairaSubscription	20000	\$20,000	30000	\$40,000	10000	\$10,000
CentralUS	20000	\$20,000	30000	\$40,000	10000	\$10,000
SairaRG	20000	\$20,000	30000	\$40,000	10000	\$10,000
Test-123	10000	\$10,000	10000	\$10,000	10000	\$10,000
Test-126	10000	\$10,000	20000	\$30,000		
Azure Monitor	10000	\$10,000	10000	\$10,000	10000	\$10,000
Virtual Machines	50000	\$50,000	50000	\$55,000	20000	\$20,000
Total	80000	\$80,000	90000	\$105,000	40000	\$40,000

Fig. 4 Azure Monthly Usage Cost by Service

5.1. Leveraging Policy Engines

Newer cloud financial management tools like CloudHealth and Cloudability include advanced policy engines that will turn their governance frameworks into common and enforceable rules. The solutions offer more than passive monitoring by taking real-time remediation steps, including termination of non-compliant resources or initiation of stakeholder alerts if a predetermined cost or usage limit is exceeded. The active enforcement makes policy enforcement doable at scale without as much human monitoring.

5.2. Implementing Guardrails and Approval Workflows

As cloud usage ramps up within enterprises, setting up active financial controls for using automated guardrails and scaffolded approval functions has shifted from best practice to business necessity. Technical Program Managers are also key agents of change in this endeavor, and they work to unify operations across functions by setting goals with cross-functional partners to:

- Instituting a budget warning and expenditure limits by team or project
- Insisting that business resources costing much money must be justified.
- Automating business hours out-of-hours shutdown of the production environment

5.3. Fostering a Culture of Cost Awareness

TPMs help tremendously in the development of cost-conscious engineering minds. Some of the practices are;

- Put cost KPIs in the OKRs of the engineering team
- The activities of cost reduction, identification and recognition

- Encouraging good sharing of knowledge regarding best practices when it comes to cost optimization

5.4. Case Study: Launch Darkly Flag Deprecation

One company ran a program to delete unused feature flags on its Launch Darkly system. By setting the goal of removing 60 percent (770) of flags, they reduced their Launch Darkly licenses to 100, cutting about 130,000 dollars per year. The other licenses had to be dealt with meticulously, and the engineering units had to collaborate effectively on this project.

6. Driving Continuous Optimization

To manage clouds, one needs to optimize, rather than make single-time fixes continuously. To ensure their cost-effectiveness when workloads change, leading organizations employ a cyclical FinOps process, including usage pattern analysis, rightsizing activities, and discount optimization. These include automated recommendations, e.g. using AWS Compute Optimizer or Kubernetes Vertical Pod Autoscaler, and regular (i.e. periodic) optimization sprints, which put teams into the active mindset of being eager to eliminate waste. With the inclusion of cost reviews in DevOps operations and the motivation of the engineering team through showback reporting, organizations develop a culture of having financial responsibility in conjunction with innovation. Such a continuous improvement strategy normally results in 20-40 percent long-term savings and avoids the cost creep resulting from new deployments.

6.1. Implementing Automated Cost Reduction Strategies

The TPMs should work with the infrastructure teams to implement automated cost-cutting measures (e.g.,

- Planning the start and finish of resources that do not work in finishing manufacturing
- Automatic unneeded volume and snapshots decommissioning

Lessen bloated database servers

6.2. Establishing a FinOps Practice

Management of the cost is not a once-in-a-lifetime event. The next thing that the PMs should pursue is the creation of a formal FinOps (Cloud Financial Operations) unit within the business. Such an interdisciplinary unit can:

- Cloud costs tracking and cloud costs trend calculations
- Identify and position those areas that need to be done better
- Connect and Mobilize budgetary programmes

Train and provide documentation to working parties of engineers

6.3. Using Cloud Provider Cost Optimization Tools

The cost optimization tools offered by the major cloud providers are natively something that TPMs should take advantage of:

- Trusted Advisor and AWS Cost Explorer

- Advisor and Azure Cost Management
- Recommender and Cost Management from Google Cloud

6.4. Case Study: S3 Bucket Lifecycle Policies

A single company has saved over 100,000 per annum through its top 10 largest buckets via S3 bucket lifecycle policies that include shifting of data not frequently implemented into lower cost storage surfaces automatically.

7. Balancing Cost and Performance

Optimization of costs in cloud environments without a sacrifice on performance is a strategic approach that addresses how clouds are structured in relation to the actual workload needs. An overprovision would guarantee performance but would use unnecessary money, and an excessive cost-cut would result in less efficient user experiences and operational failures.

The strategies that work efficiently are setting up auto-scale policies to add and subtract resources dynamically, choosing cost-effective instance types depending on the pattern of workloads and using performance-monitoring tools to find instances of inefficiency. By performing analysis of tradeoffs such as load testing with various resources and using tiered storage solutions, organizations can strike the right balance of maintaining responsiveness and yet doing away with unnecessary expenses. This tradeoff is especially critical in applications facing external customers where the performance is a direct determinant of revenue, and continuous cost-performance analysis of cloud systems is a crucial skill set in cloud teams.

7.1. Utilizing Observability Tools

Cost data can then be correlated with performance indicators via such platforms as New Relic and Datadog and used to make data-driven capacity planning decisions. This plan ensures that cost tradeoffs are in line with the company's corporate goals and SLOs.

7.2. Defining the Right Metrics

In as much as cost management is important, it cannot be pursued at the cost of performance and dependability. The PMS must cooperate with SRE teams to establish and observe metrics that will balance the cost and performance. Some of the key considerations to be made are:

- Error budgets and service level goals (SLOs)
- Cost per API call or transaction
- Resource utilization vs. response time curves

7.3. Putting Performance-Aware Autoscaling into Practice

Over-provisioning may happen when forms of traditional autoscaling are only dependent on either CPU or Memory usage. TPMs must encourage more advanced autoscaling algorithms with consideration of application-specific increment/decrement measures and historical trends.

7.4. Case Study: Machine Learning Model Optimization

A company expects to save a million times 800,000 yearly, through better organization of machine learning. The more important strategies were

- Optimizing the number of data loading to ML training
- When possible, go for a smaller instance size
- It involves undertaking an automated type of shutdown of resources during off periods.

This initiative proved that optimization of costs and effectiveness of operation complement each other by reducing costs and increasing the effectiveness of training processes.

8. Cloud Migration and Modernization

Moving to the cloud does not imply merely transferring and extending the current infrastructure but presents a chance of renovating applications and reducing costs through cloud-native designs. The migration process needs to be gradual: to make the shift between worlds, first, determine the compatibility of workloads, refactor an application into microservices, and utilize the platform, such as managed databases and serverless, to eliminate the burden of operations. Rehosting (lift-and-shift) is a good way to achieve quick returns. However, the real transformation, where all the benefits of elasticity, automation and pay-as-you-go contracts are leveraged, is achieved by re-platforming and re-architecting. Modernization also unleashes advantages such as enhanced scalability, accelerated deployment rates, and efficient resource use, although it needs to be planned carefully so that both the short-term expenditure and the long-term ROI are taken into account. This means that by structuring a migration strategy according to business priorities, whether focused on speed, cost reduction, or innovativeness, organization can make the most out of cloud value and the least out of operation interference.

8.1. Utilizing Cloud-Native Services

Teams should aim to leverage the deployment of cloud-native services as they often prove to be more cost-effective than the so-called lift-and-shift approaches, which should be encouraged to be employed by TPMs. These are some of the examples:

- A changeover to the management of databases as opposed to self-managed databases
- Possible applications of server-free computing to on-demand jobs
- To improve the use of resources using container orchestration platforms

8.2. Calculating the Overall Cost of Ownership (TCO)

When moving from on-premises or legacy cloud deployments, TPMs must perform a detailed TCO analysis. TCO needs to be independently done by TPMs when trying to move from on-premise or legacy cloud environments. It ought to contain the following:

- Immediate costs of infrastructure
- Charges for licenses and maintenance
- Running costs (e.g., personnel, training)
- Expenses related to migration and possible downtime

8.3. Case Study: ROSA to EKS Migration

A substantial migration to one organization was initiated with Red Hat OpenShift Service on AWS (ROSA).

Amazon Elastic Kubernetes Service (EKS). These were some of the main discoveries:

- Elimination of the licensing fees of ROSA (about 900,000 dollars per annum)
- Optimized usage of resources with the help of the scaling capabilities of Kubernetes
- Reduce operational expense by using the AWS control plane

During the first part of the project, 60 percent of deployments were transferred, and 84 percent of services are working in production in EKS. This demonstrates why making smart decisions and moving to cloud services can lead to substantial cost savings.

9. Vendor Management and Contract Optimization

Cost negotiation is not enough to achieve effective cloud vendor management, as the establishment of contracts has to follow the usage pattern of the cloud along with business and technical needs, which keep changing. To remain viable and competitive, organizations need to constantly review consumption patterns so that they maximize such commitments as Reserved Instances, Savings Plans and enterprise discounts without over-committing.

Formal vendor review, such as performance comparison, support service level agreements and region capabilities, prevents lock-in and creates a balance of leverage during negotiations. Leveraging FinOps best practices and approaches through contract reviews, multi-cloud hedging practices, and granular chargeback reporting, companies can optimize cloud expenditures by up to 30-40 percent, and the contractual terms are much more likely to satisfy the varied business requirements at the current and longitudinal capacities. The best teams do this all the way up to automatic monitoring of the commitment, with tools such as ProsperOps or vendor-specific advisors making real-time changes to purchases as usage changes.

9.1. Software Licensing Optimization

The terms of software licenses for tools that integrate with their cloud services should fall within regular reviews by TPMs. As examples of techniques, there are:

- Having projects or teams under one roof concerning licenses
- The investigation of the open source equivalents, whenever possible
- Negotiating pricing models on the basis of use

9.2. Utilizing Enterprise Discount Programs

The negotiation of the Enterprise Discount Programs (EDPs) can lead to high savings for enterprises that spend

substantially on cloud services. To: TPMs must work hand in hand with procurement teams in working towards the following:

- Analyze the use of cloud in the present and the future.
- Get a cut-off of opportunities in volumes.
- Negotiate flexible terms that consider the curves, boom and bust and growth.

9.3. Case Study: Multi-Vendor Cost Optimization

Vendor management: An organization established a comprehensive vendor management plan, which resulted in significant cost savings.

- Datalog: Solved a \$1.1 million invoicing matter and came to a new agreement.
- LaunchDarkly: \$130,000 saved each year in the number of licenses
- Migration of all FullStory data to Datalog Session Replay: Enhanced feature set and \$150,000 of annual savings

It is evidence of the reward of being assertive in terms of handling vendor management and predictions of driving down the costs by appropriate instrument choice and contract negotiation.

10. Building a Cost-Aware Engineering Culture

The first milestone toward building a cost-sensitive engineering culture is understanding cloud usage at all times via dashboard, showback analytics, and per-project, per-feature cost allocation. It requires engineers to remember the cost metrics at all phases of development, beginning with architecture reviews and extending through deployment pipelines, so that engineers can make tradeoffs based on knowledge of performance, scalability, and efficiency.

Accountability is set by gamification, including those cost-saving challenges or peer recognition based on optimization victories. In contrast, FinOps-foundation training can help teams learn how to identify wastage (e.g. decommission idle resources and oversized instances). Leadership is essential as it has to encourage both cost efficiency and innovations, making it a common KPI instead of a side note.

In the long term, this cultural change will turn cost optimization into an implicit part of the engineering process, thus turning it into a source of sustainable savings that do not undermine agility.

10.1. Integrating Cost Awareness into the Development Lifecycle

The cost factor should be considered by TPMs in all areas of the software development process:

- Adding the projections of costs to the feature planning and prioritization
- Checking costs by running the CI/CD pipelines
- Carrying out regular cost retrospectives in addition to regular sprint retrospectives

10.2. Education and Training

TPMs help educate the teams about learning engineering and the financial implications of their technical decisions. Strategies include:

- Organizations of training on the principles of cloud economics
- Sharing of the successful experience of cost savings
- Offering the guidelines that are cost-effective for architectural design and development processes

10.3. Gamification and Incentives

It is possible to introduce game mechanics to enhance the involvement in cost optimization projects:

- Rewards or acknowledgement of cost-saving creative thinking
- Leaderboards of the most cost-efficient teams
- Cost Optimization Hackathons as a method to discover new opportunities

10.4. Case Study: Engineering OKRs for Cost Optimization

One organization achieved this by doing the following:

- Incorporation of team-level goals of cost optimization in OKRs
- Provision of live cost dashboards to each engineer
- The optimization of costs triumphed in the organizational forums through celebration and the spread of news

Optimization spearheaded by engineers contributed to the primary factor influencing the twenty percent reduction of cloud costs in six months.

11. Future Trends and Research Directions

Artificial intelligence-led automation, sustainable computing and distributed computing will define cloud financial management in the future. Innovative technologies such as predictive autoscaling (based on ML and designed to deliver precise forecasts on demand patterns) and intelligent workload placement (seeking the equilibrium between cost-performance ratios and reduced eco-friendly output) could revolutionize cost optimization.

There is a growing interest in real-time FinOps, where autonomous systems will automatically provision resources and make commitments, and blockchain-based cloud cost transparency tools will allow cross-organizational and audit-able spending benchmarks. In the meantime, the emergence of serverless orchestration and microservices, with millisecond billing, will necessitate new models to track fine-grained cost tracking.

Industry and academic partnerships If you compare the current state with what is being explored by organizations like the FinOps Foundation through their research initiatives, one of the most apparent changes across the industry will be an approach towards how quantum

computing and edge-native architecture will prolong the survival of cloud economics in its current form, as in this case, continuous innovation will take place to keep up with how the field of infrastructure development is being transformed.

11.1. Sustainability and Cost Optimization

With companies showing an increasing concern over the environmental impact they have, studies being developed need to consider the relationship between cost reduction and sustainability in the context of cloud computing:

- Quantifying the amount of carbon emitted by cloud tasks
- The creation of optimization plans, which consider both cost, eco-friendliness and performance, is an essential redesign task.
- The field of research on the effect of renewable energy on cloud pricing models is known as research on the impact of renewable energy on cloud pricing models.

11.2. AI-Driven Cost Optimization

With the advancement of technology in the field of artificial intelligence and machine learning, more sophisticated and automated techniques of optimizing cloud costs are bound to come into practice. Future research is required to look into:

- Predictive analytics used to have some insight into the allocation of resources in the future
- Artificial intelligence as a workload-distribution system across a number of cloud platforms
- Adaptation of applications (cost-efficiently) automatically

11.3. FinOps Maturity Models

What is needed is consistency in maturity models aligned to cloud financial operations. Investigations in this area may be focused on:

- Outlining the essential skills and measures at every level of maturity
- Developing assessment tools for companies to benchmark the FinOps approaches
- Realizing good strategies to be used to progress along the maturity stages

12. Comparative Analysis

Emerging findings on cloud cost optimization have been published during recent years, but tend to lack integration since they are more narrowly specialized in the aspects of autoscaling, predictive provisioning or resource scheduling rather than presenting them as a single coherent framework, as in the case of SISSGECO.

The following is a more literature-informed, comparative account of the same:

Table 1. Comparative Study – Existing Methods vs SISSGECO

Authors	Method	Limitations	SISSGECO Value
Aslanpour et al. (2018)	Cost-aware scaling via MAPE loop	Rule-based; only on execution	ML-driven across the full MAPE cycle
Pan et al. (2023)	GP-based autoscaler (Alibaba ECS)	Vendor-specific; no feature selection	Multi-cloud + hybrid ML + features
FGCS (2018)	Microservice resource sizing	Shallow classifiers; weak features	Enhanced selection + multi-model
Ghasemi et al. (2023)	Learning automata for provisioning	General; not workload-focused	Real-time, workload-based tuning

12.1. Novelty of the Work

- Comes up with a universal model that works over multi-cloud and hybrid infrastructure, which does not lock in to vendors.
- Embeds make optimization of cost a primary goal, not just one of the metrics to monitor.
- Applies Improved ReliefF and Mutual Information-based Feature Selection to transfer the workload in such a way that it is exactly mapped as a known set of cost-performance values.
- Aixedes mix and match a set of neural networks (NN), random forest (RF), recurrent neural networks (RNN), and K-NN to achieve intelligent classification by optimally assigning weights.

13. Conclusion

With the costs of cloud services rapidly growing as a percentage of total IT budgets, Technical Project Managers (TPMs) now have an excellent opportunity to help their companies deliver more business value by

effectively managing the costs of such services. With technical flair combined with financial expertise, TPMs would help companies strike the right mix of innovation and cost-efficient enterprise in the cloud era.

Individuals who become skilled in cost transparency, resource efficiency, and financial oversight will be in good standing as the need for cloud financial management increases. The best TPMs would be those that can fill in the divide between the technical and financial conversation by translating computer data and financial numbers.

The techniques and illustrations presented in this article can provide a good beginning for TPMs who want to lead cost reduction initiatives. This area of business is evolving fast; however, constant education and exchange of experiences will be necessary to keep up with the trends. The TPM's role in cost management will also grow as developments in cloud technology and financial strategies evolve.

References

- [1] Davide Taibi, and Kari Systä, "From Monolithic Systems to Microservices: A Decomposition Framework based on Process Mining," *International Conference on Cloud Computing and Services Science*, pp. 153-164, 2019. [[Google Scholar](#)] [[Publisher Link](#)]
- [2] Giovanni Toffetti et al., "Self-managing Cloud-native Applications: Design, Implementation, and Experience," *Future Generation Computer Systems*, vol. 72, pp. 165-179, 2017. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [3] Yazn Alshamaila, Savvas Papagiannidis, and Feng Li, "Cloud Computing Adoption by SMEs in the North East of England: A Multi-perspective Framework," *Journal of Enterprise Information Management*, vol. 26, no. 3, pp. 250-275, 2013. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [4] Armin Balalaie, Abbas Heydarnoori, and Pooyan Jamshidi, "Microservices Architecture Enables DevOps: Migration to a Cloud Native Architecture," *IEEE Software*, vol. 33, no. 3, pp. 42-52, 2016. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [5] Trevor Clohessy, Thomas Acton, and Lorraine Morgan, "The Impact of Cloud-based Digital Transformation on IT Service Providers: Evidence from Focus Groups," *International Journal of Cloud Applications and Computing*, vol. 7, no. 4, 2017. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [6] Tomislav Vresk, and Igor Cavrak, "Architecture of an In-teroperable IoT Platform based on Microservices," *2016 39th International Convention on Information and Communication Technology, Electronics and Microelectronics*, 2016. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [7] Jacopo Soldani, Damian Andrew Tamburri, and Willem-Jan Van Den Heuvel, "The Pains and Gains of Microservices: A Systematic Grey Literature Review," *Journal of Systems and Soft-ware*, vol. 146, pp. 215-232, 2018. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [8] Yevheniya Nosyk, "Migration of a Legacy Web Application to the Cloud," Theseus, 2018. [[Google Scholar](#)] [[Publisher Link](#)]
- [9] Ahmed Barnawi et al., "The Views, Measurements and Challenges of Elasticity in the Cloud: A Review," *Computer Communications*, vol. 154, pp. 111-117, 2020. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [10] Davide Taibi, Valentina Lenarduzzi, and Claus Pahl, "Processes, Motivations, and Issues for Migrating to Microservices Architectures: An Empirical Investigation," *IEEE Cloud Computing*, vol. 4, no. 5, pp. 22-32, 2017. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]

- [11] Sai Dikshit Pasham, "AI-Driven Cloud Cost Optimization for Small and Medium Enterprises (SMEs)," *The Computertech*, vol. 3, no. 1, 2017. [[Google Scholar](#)] [[Publisher Link](#)]
- [12] FinOps: A New Approach to Cloud Financial Management, O'Reilly Media, 2020. [Online]. Available: <https://cxo-institute.com/wp-content/uploads/2020/07/FinOps-A-New-Approach-to-Cloud-Financial-Management-1.pdf>
- [13] Lee Atchison, *Architecting for Scale: High Availability for Your Growing Applications*, O'Reilly Media, 2016. [[Google Scholar](#)] [[Publisher Link](#)]
- [14] Martin Fowler, *Patterns of Enterprise Application Architecture*, Addison-Wesley Professional, 2002. [[Google Scholar](#)] [[Publisher Link](#)]
- [15] Jez Humble, and David Farley, *Continuous Delivery: Reliable Software Releases through Build, Test, and Deployment Automation*, Addison-Wesley Professional, 2010. [[Google Scholar](#)] [[Publisher Link](#)]
- [16] Gene Kim et al., *The DevOps Handbook: How to Create World-Class Agility, Reliability, and Security in Technology Organizations*, IT Revolution Press, 2016. [[Google Scholar](#)] [[Publisher Link](#)]
- [17] The Twelve-Factor App. [Online]. Available: <https://12factor.net/>
- [18] Naga Surya Teja Thallam, "Comparative Analysis of Public Cloud Providers for Big Data Analytics: AWS, Azure, and Google Cloud," *International Journal of AI, Bigdata, Computational and Management Studies*, vol. 4, no. 3, pp. 18-29, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]